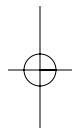


VOLUME 76, NUMBER 1, JULY 2009



ISSN 0887-3585

Articles published online in Wiley InterScience, 21 October 2008–24 February 2009

Conserved amino acid networks involved in antibody variable domain interactions

Norman Wang,[†] William F. Smith,[†] Brian R. Miller, Dikran Aivazian, Alexey A. Lugovskoy, Mitchell E. Reff, Scott M. Glaser, Lisa J. Croner,* and Stephen J. Demarest*

Biogen Idec, San Diego, California

ABSTRACT

Engineered antibodies are a large and growing class of protein therapeutics comprising both marketed products and many molecules in clinical trials in various disease indications. We investigated naturally conserved networks of amino acids that support antibody V_H and V_L function, with the goal of generating information to assist in the engineering of robust antibody or antibody-like therapeutics. We generated a large and diverse sequence alignment of V-class Ig-folds, of which V_H and V_L domains are family members. To identify conserved amino acid networks, covariations between residues at all possible position pairs were quantified as correlation coefficients (ϕ -values). We provide rosters of the key conserved amino acid pairs in antibody V_H and V_L domains, for reference and use by the antibody research community. The majority of the most strongly conserved amino acid pairs in V_H and V_L are at or adjacent to the V_H - V_L interface suggesting that the ability to heterodimerize is a constraining feature of antibody evolution. For the V_H domain, but not the V_L domain, residue pairs at the variable-constant domain interface (V_H - C_H1 interface) are also strongly conserved. The same network of conserved V_H positions involved in interactions with both the V_L and C_H1 domains is found in camelid V_{HH} domains, which have evolved to lack interactions with V_L and C_H1 domains in their mature structures; however, the amino acids at these positions are different, reflecting their different function. Overall, the data describe naturally occurring amino acid networks in antibody Fv regions that can be referenced when designing antibodies or antibody-like fragments with the goal of improving their biophysical properties.

Proteins 2009; 76:99–114.
© 2008 Wiley-Liss, Inc.

Key words: immunoglobulin variable domain; Ig-fold; V-class; covariation; antibody engineering.

INTRODUCTION

Antibodies are useful targeted therapeutics because of their ability to bind specific ligands with high affinity and specificity. Antibody variable domains (V_H in the heavy chain, V_L in the light chain), which provide the binding capability, may be purposely engineered to impart desired antigen recognition or binding affinity properties. Some designs have implemented recombinant production of isolated V_H - V_L domains (Fv region), providing researchers with more design flexibility than standard antibody therapeutics (e.g., the expression of the Fv region as a single polypeptide chain or “scFv”^{1–5}). However, removal of the V_H - V_L domains from the quaternary structure of an antibody can lead to stability and solubility problems. Several mechanisms have been proposed to account for the generally poor biophysical behavior of scFvs and related designs, and include the intrinsic instability of the isolated domains, the weak affinity between V_H and V_L domains, and the absence of possibly stabilizing interactions with the antibody constant domains.⁶ An understanding of the specific amino acids that mediate interactions between the V_H and V_L domains, and between the variable and constant domains, would enable improved designs of antibodies and antibody-like proteins.

Antibody variable domains are part of the immunoglobulin domain or “Ig-fold” superfamily. The Ig-fold superfamily is a large group of structurally related protein domains commonly found in mammalian cell surface proteins or in soluble extracellular signaling proteins.⁷ Ig-fold domains consist of two β -sheets, each arranged in a “Greek-key” topology, that are packed tightly against one another and are generally supported by an intradomain disulfide bond. Depending on the number of strands in each β -sheet and the loop connections between the strands, the superfamily can be divided into several subfamilies including the C-, I-, and V-classes.⁸ Antibody variable domains are V-class Ig-folds and their constant domains are C-class Ig-folds (Fig. 1). Ig-fold or Ig-fold-like domains are also present in cell adhesion proteins, integrins, allergens, T-cell

Additional Supporting Information may be found in the online version of this article.

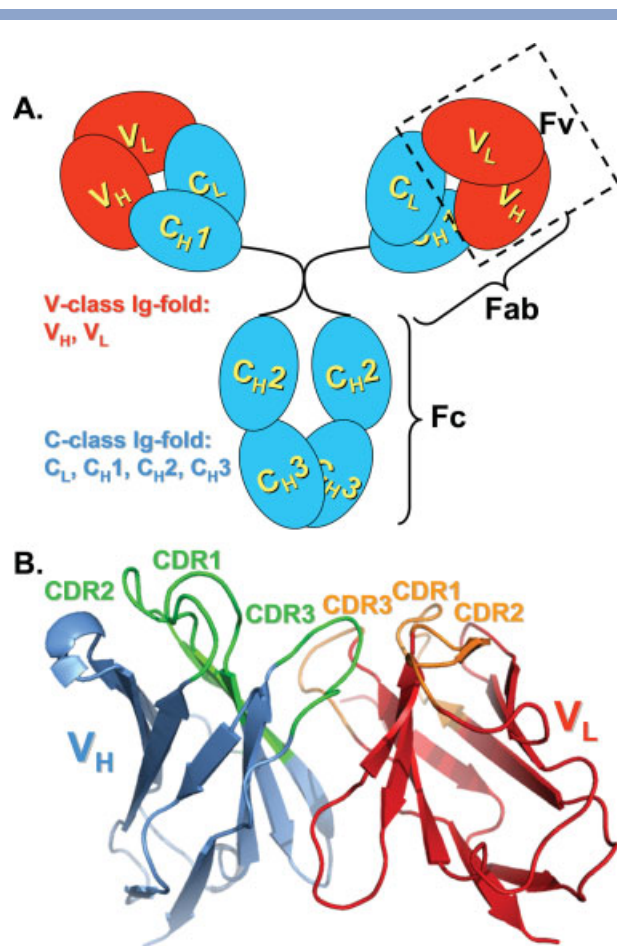
[†]Norman Wang and William F. Smith contributed equally to this work.

*Correspondence to: Stephen Demarest and Lisa Croner, Biogen Idec, 5200 Research Place, San Diego, CA 92122. E-mail: stephen.demarest@biogenidec.com and lisa.croner@biogenidec.com.

Received 5 May 2008; Revised 3 October 2008; Accepted 16 October 2008

Published online 11 November 2008 in Wiley InterScience (www.interscience.wiley.com).

DOI: 10.1002/prot.22319

**Figure 1**

Diagrams of an immunoglobulin and its Fv domain. **A:** Schematic diagram of an IgG antibody. The variable domains that compose the antigen-binding or Fv-region are shown in red and the constant domains are shown in blue. The variable domains are V-class Ig-folds, whereas the constant domains are C-class Ig-folds, which are highly similar to V-class Ig-folds, but lack two additional β -strands commonly found in V-class structures. **B:** Ribbon diagram of an antibody Fv-region consisting of a variable domain from the immunoglobulin heavy chain (V_H-blue) and a variable domain from the immunoglobulin light chain (V_L-red). The complementarity determining regions (CDRs) of the V_H (shown in green) and the V_L (shown in orange) comprise the antigen-binding site.

receptors, major histocompatibility complexes, immunoglobulin receptors, and many other protein families with diverse functions.

The past decade has seen a significant increase in the number of publicly available Ig-fold sequences. Large databases of antibody variable domain and T-cell hyper-variable domain Ig-fold sequences have been compiled.⁹ Information from these databases has been instrumental in antibody humanization, affinity maturation, and the stabilization of single chain Fv (scFv) or other antibody constructs.^{5,10–13} Antibody sequence databases generally influence antibody design by enabling frequency analyses

at single amino acid positions (i.e., consensus modeling) that may be used for generating rational designs.^{12,14–16} Recent studies with other protein domain superfamilies have extended sequence-based approaches by examining how amino acid *pairs* or *networks* may be conserved within subsets of a protein superfamily with related function or across diverse members of a protein superfamily. Such amino acid networks may define important structural or functional features of these protein domains.^{17–19} These approaches, sometimes referred to as “covariation analyses,” track whether the presence (or absence) of a particular amino acid at one position correlates with the presence (or absence) of another amino acid at a second position within a multiple sequence alignment. Although covariation analyses have been performed on several protein families (including SH3 domains, WW-domains, TPR-motifs, GPCRs, serine proteases, globins, viral coat proteins, and others^{20–23}), very little has been described concerning covariation analyses of Ig-folds. The paucity of Ig-fold covariation data may stem from several factors, one being that large collections of Ig-fold sequences were, until recently, limited primarily to antibody sequences, particularly human and murine.^{24,25} Also, accurate alignment of diverse members of large proteins (>100 amino acids) like Ig-fold domains is challenging and misalignments can limit the validity of covariation data.²³

Here we describe the application of covariation analyses to a high-quality, 3D-structure-based alignment of diverse V-class Ig-fold sequences. A diverse V-class Ig-fold sequence alignment was constructed, and covariations were quantified as correlation coefficients (ϕ -values²³) for every amino acid pair (i.e., every residue combination found at all possible pairs of positions) in the alignment. The results serve as a rich repository of amino acid interactions conserved throughout Ig-fold evolution. The data reveal conserved residue networks that may support interactions between the V_H and V_L domains. The data also reveal that V_H domain networks involved in interactions with V_L domains are co-conserved with residue networks observed at the V_H–C_H1 junction suggesting that these two functional areas have coevolved to support the overall quaternary antibody structure.

METHODS

Creation of structure-based Ig-fold alignments

Structures of Ig-fold proteins or Ig-fold domains from multidomain proteins were gathered from the ASTRAL database,^{26,27} which contains domain structures matching the Structural Classification of Proteins (SCOP, Version 1.69) hierarchy. SCOP classifies the Ig-fold as a member of the “Beta Proteins, Immunoglobulin-like β -sandwich fold, Immunoglobulins” superfamily.²⁸ PDB

files of the V-class Ig-folds were downloaded using customized shell scripts. Each Ig-fold structure was inspected visually using Swissprot DeepView; sequences were removed from the study if they were erroneously categorized, incomplete (either missing residues due to a lack of electron density or domain swapped²⁹), redundant (i.e., those with identical sequences), or obviously did not conform to the β -sandwich Ig-fold topology. Sequences of aberrant length (>2 -times the standard deviation about the mean V-class length, 112.0 ± 10.6 residues) were also removed. 702 structures were aligned using the Secondary Structure Matching (SSM) assisted implementation in the Schrödinger Prime structalign program.^{30–32} Schrödinger Prime was used to generate structure-based V-class sequence alignments based on the proximity of each C_{α} atom subsequent to an all-to-all structure alignment, which minimized the average distance between all structural pairs. The alignment was most accurate in the regular β -strand regions and less accurate in the connecting loop regions because of variable loop lengths and structures.

Generation of a diverse V-class Ig-fold sequence alignment

A custom Hidden Markov Model (HMM) of the V-class Ig-folds was built from the structure-based sequence alignments. The HMM was created with the HMMER software package (version 1.8), using the “hmmbuild” and “hmmcalibrate” functions.³³ The HMM was used to find potential V-class Ig-fold sequences in the NR-database maintained at NCBI using the “hmmsearch” function. The output of this function ranked the hit sequences by their scores relative to our custom V-class HMM, and sequences with scores above a recommended threshold were retained as candidate members of the V-class dataset. The output also provided the number of “hits” per sequence (i.e., the number of Ig-folds within a contiguous gene sequence) and the exact residue positions of the hits. For sequences containing one or more candidate V-class sequences, the relevant subsequences were extracted from the full NR sequence using a custom Java executable. As an additional test to confirm that each sequence pulled from NR using our V-class HMMs belonged to the V-class Ig-fold subfamily, the custom shell script “pfamverify” using the HMM tool “hmmpfam” was applied to each Ig-fold candidate sequence.³⁴ Ig-clan HMMs (including V-, I-, C1-, C2-, and less specific Ig HMMs) were downloaded from the PFAM website. Sequences that scored lower with the PFAM V-class HMM than with other PFAM Ig-fold HMMs, and sequences whose score with the PFAM V-class HMM lay below recommended cutoffs (TC1 defined at the PFAM website) were removed. Thus, V-class Ig-fold sequences were retained only if their PFAM scores validated their Ig-fold class assignments. The Ig-

fold sequences extracted from NR were aligned by our structure-based V-class HMM using the “hmmalign” function in the HMMER package. The resulting V-class dataset contained 48,696 sequences including those from both the SCOP 3D protein database and NR.

The resulting sequence collection was biased toward well-studied Ig-fold-containing proteins (i.e., human and murine V-class sequences frequently deposited in NR). To reduce the over-representation of these sequences, we developed a heuristic algorithm that eliminated sequences based on identity cut-off criteria. In brief, percent identities were calculated for all sequence pairs. Sequences were grouped into bins representing their maximum percent identity with any other sequence (i.e., 99% bin, 98% bin, 97% bin, etc.). Sequences within each bin were then ranked according to decreasing nongap residue count, giving better ranks to sequences with fewer gaps. In each bin, sequences with an equal number of nongap residues were ranked by Henikoff weights to filter out more common sequence types while preserving rare sequences with the goal of increasing diversity within the final datasets.³⁵ An identity cutoff of 80% was used for V-class sequences. This left 2786 sequences, each with less than 80% identity to all other sequences, in the V-class dataset.

The resulting multiple sequence alignment contained many positions populated by gaps ($>50\%$ gaps for most sequences). To eliminate this problem, columns that were not match states in the HMM were removed. This resulted in 144 remaining columns for our custom V-class alignment. Still, 354 sequences contained $>40\%$ gaps. These sequences, which were generally incomplete, were removed from the alignment. The final V-class Ig-fold dataset contained 2432 sequences. Virtually all the V-class sequences in the dataset were naturally occurring (nonengineered).

Correlation coefficient (ϕ -value) calculation

Covariation between amino acid pairs in multiple sequence alignments were calculated as correlation coefficients (ϕ -values), as described previously.²³ The calculations were encoded into a Java executable and run with Java Runtime Engine (JRE) version 1.4.2. ϕ -values were defined as

$$\phi(x_i y_j) = \frac{(x_i y_j \times \bar{x}_i \bar{y}_j) - (x_i \bar{y}_j \times \bar{x}_i y_j)}{\sqrt{(x_i y_j + \bar{x}_i \bar{y}_j) \times (x_i \bar{y}_j + \bar{x}_i y_j) \times (x_i y_j + x_i \bar{y}_j) \times (\bar{x}_i y_j + \bar{x}_i \bar{y}_j)}}, \quad (1)$$

where $x_i y_j$ is the number of times amino acids of type “ x ” or “ y ” are found in the same sequence at positions i and j , respectively, $\bar{x}_i \bar{y}_j$ is the number of times both amino acids are absent from the same sequence, $x_i \bar{y}_j$ is the number of times x is found present while y is absent,

and $\bar{x}_i y_j$ is the number of times x is absent while y is present. This equation can be rewritten as:

$$\phi(x_i y_j) = \frac{(a \times d) - (b \times c)}{\sqrt{efgh}}, \quad (2)$$

where a to h are given by the contingency table:

	x_i	\bar{x}_i	Total
y_j	a	b	e
\bar{y}_j	c	d	f
Total	g	h	

and $a = x_i y_j$, $b = \bar{x}_i y_j$, $c = x_i \bar{y}_j$, $d = \bar{x}_i \bar{y}_j$, $e = a + b$, $f = c + d$, $g = a + c$, and $h = b + d$. Particular residue pairs (specific combinations of residues at specific positions) were not considered unless they were observed in the alignments a minimum of 10 times.

Statistics

Statistical significance of the ϕ -values was evaluated with a chi-square (χ^2) test, using Bonferroni-corrected P -values to adjust for multiple testing.³⁶

The χ^2 test is often used to evaluate the significance of values observed in contingency tables of two dichotomous variables, such as the contingency table above. The equation for this use of χ^2 can be written as

$$\chi^2 = \sum \left[\frac{(o_k - e_k)^2}{e_k} \right], \quad (3)$$

where o_k stands for the observed frequency and e_k stands for the expected frequency in one cell of the table. χ^2 is calculated by taking the sum of the squared and normalized differences between the observed and expected frequencies over all the cells. When expected frequencies are unknown, they can be estimated from observed frequencies and the equation becomes

$$\chi^2 = \frac{N \times (ad - bc)^2}{efgh}, \quad (4)$$

with a to h representing the same values as in Eq. (2) above, and N standing for the total number of samples. Comparing Eqs. (2) and (4), it is evident that

$$\chi^2(x_i y_j) = \phi(x_i y_j)^2 \times N$$

This relationship between χ^2 and ϕ is useful because of the rich information available about the χ^2 statistic, including tables of P -values for χ^2 with specified degrees of freedom (df). However, before proceeding to use this

relationship, we performed random simulations to confirm that it held for our dataset. We took our V-class sequence alignment, performed repeated (tens of thousands) random shuffling of the residues at each position (so that the residue frequencies at each position remained unchanged, but the correlations across positions were randomized), and computed ϕ -values for each randomization. We then calculated the probabilities of observing specific strong covariations by chance directly from these random simulations. In all cases examined, we found close agreement between the probabilities observed in the simulations and those calculated from χ^2 .

Having validated the use of χ^2 to determine significance in our dataset, we converted our ϕ -values to χ^2 using Eq. (4), and used standard χ^2 tables (df = 1) to find P -values. ϕ -values were calculated for the 186,171 amino acid pairwise combinations that occurred at least 10 times in our V-class alignments. To correct for multiple testing, we used a Bonferroni-corrected P -value³⁶ as our criterion for significance, striving for significance at the true $P < 0.0001$ level. For positive correlations, this corresponded to ϕ -values > 0.1237 . Use of the Bonferroni correction along with this strict P criterion gave a very conservative list of amino acid pairs with significant ϕ -values, and greatly reduced our chances of finding false positives. Even with this conservative approach, 13,796 significant positive correlations (see Results) were observed.

RESULTS

Alignment quality and diversity

Information about protein 3D structures can significantly improve the quality of multiple sequence alignments.³⁷ As described in the Methods, we compiled 3D structures of V-class Ig-folds from SCOP and generated a structure-based multiple sequence alignment of these V-class domain sequences. A custom HMM was built from this structure-based alignment and used to align additional Ig-fold sequences from the NR sequence database. Here we discuss the quality and the diversity of the resulting alignment.

The quality of the alignment guided by our structure-based HMM was evaluated by examining whether the HMM properly aligned sequences that, though disparate in residue identity, are known to form the same part of an Ig-fold 3D structure. As expected, residues that make up the β -strands of both V_H and V_L domains were well aligned, whereas alignments of residues in the loop regions, whose structures are more variable, often contained many gaps. We also looked to see whether the HMM properly aligned the consensus sequences of antibody V_H and V_L chains, in the V_H – V_L interface region. The heterodimeric structure of the interface is highly symmetrical (both V_H and V_L domains use the same face

of their Ig-fold to create the heterodimer). Thus, residues buried in the interface should align in 3D, despite differences between the amino acids of V_H and V_L Ig-folds at these positions. We used a 1.8 Å crystal structure of an antibody Fab from our lab (Jordan *et al.*, manuscript in preparation) to determine the residues of both V_H and V_L that bury surface area at the interface using the program MOLMOL.³⁸ The V_H and V_L residue positions buried at the interface mapped to the same residue positions within the multiple sequence alignment, even though the amino acid identities at these positions vary.

Covariation analyses are most successful when applied to sequence datasets that are highly diverse.^{23,39–41} From a practical standpoint, ϕ -value correlation coefficients increase when there are many instances of both the presence and absence of a conserved amino acid pair across a sequence alignment (see numerators of Eqs. 1 and 2). Highly diverse sequence sets are more likely to contain sequences both with and without the pair than are datasets of highly related sequences. As our goal was to investigate conserved amino acid networks in V_H and V_L domains, it was therefore important that our V-class Ig-fold dataset contains a background of Ig-fold family members that are not evolutionarily constrained to perform the same function as V_H and V_L domains. Additionally, covariation signals pertaining to polar interactions have been shown to be stronger in datasets with moderate to high evolutionary distances between sequences.⁴⁰ Protein interfaces, in which Ig-folds frequently appear, often rely more heavily on polar interactions than do protein cores,⁴² providing another reason for generating a V-class Ig-fold sequence dataset that was highly diverse. The unfiltered sequence set from NCBI contained ~50,000 sequences highly biased toward immunoglobulin variable domains. An 80% identity filter was used to reduce bias toward over-represented V-class families, thus promoting diversity. After filtering, the dataset contained 2432 V-class sequences. Members of the V-class dataset could be divided into three functional categories: (1) 50% were immunoglobulin variable genes (including both V_H and V_L antibody domains); (2) 16% were T-Cell Receptor V-class genes; and (3) 34% were V-class genes derived from diverse functional families, each of which comprised less than 5% of the V-class sequences. The sequences were derived from species ranging from cartilaginous fish to primates. There was a bias toward human immunoglobulin variable domain sequences with 574 of the total 2432 V-class sequences being human V_H (484 of 993 V_H) or human V_L (90 of 187 V_L). Examples of other species contributing V_H and V_L sequences to the V-class database include mouse (44), cow (16), camel (174), llama (83), macaque (17), and chicken (9). Despite the bias toward human V_H and V_L , the average sequence identities within V_H and V_L subgroups were low – $41 \pm 1\%$ and $29 \pm 1\%$, respectively. The distribution of germline V_H and V_L sequences pass-

ing the 80% identity filter roughly matched the naturally observed distribution of variable domain sequences⁴³ suggesting that variable gene subclasses were similarly diverse and fairly represented in the sequence dataset. Positional entropy calculations using the final V-class dataset demonstrate the much higher positional diversity with the V-class dataset compared with antibody V_H or V_L datasets that may be used for consensus analyses^{15,44} (Supp. Info. Fig. 1).

Correlation coefficients (ϕ -values) between residues of V-class Ig-folds

We adopted a previously described method, the use of ϕ -value correlation coefficients, for quantifying covariations of residue pairs within sequence alignments.²³ Figure 2 shows the number and distribution of ϕ -values calculated for amino acid pairs that were observed ≥ 10 times within the V-class Ig-fold alignment. Positive ϕ -values represent positive correlations (the presence of one amino acid at one site in the alignment is correlated with the presence of another amino acid at another site). As ϕ -values move from 0 to +1, the strength of the correlation between the two amino acids increases. Negative ϕ -values represent negative correlations (the presence of one amino acid at one site in the alignment is correlated with the *absence* of another amino acid at a second site), which become stronger as ϕ -values move toward -1.0. Our statistical analyses (see Methods) showed that ϕ -values greater than 0.1237 were significant (P -values < 0.0001). This conservative statistical estimate revealed 13,796 significant positive covariations. However, statistical significance does not entirely indicate the strength of

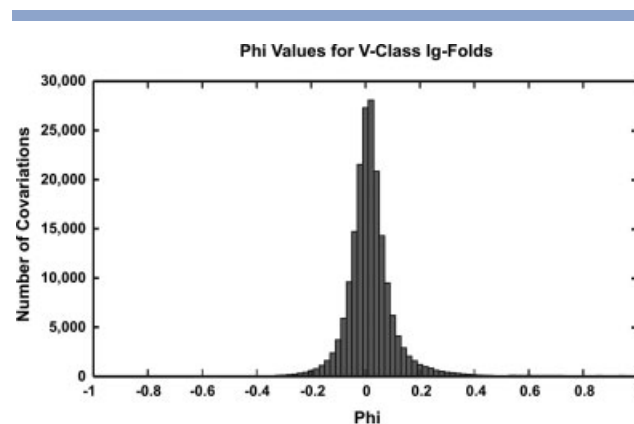


Figure 2

Distribution of ϕ -values calculated for the V-class alignment. There are 4,118,400 ($20 \times 20 \times \sum_{n=1}^{144} n - 1$) possible amino acid pairings within the V-class sequences. Of these possible pairings, 1,098,890 actually exist within the sequence database (i.e., some amino acids pairings are not observed across columns of the alignment). The histogram shows the distribution of ϕ -values from the 186,171 pairings that occur at least 10 times. The 13,796 ϕ -values greater than 0.1237 were considered statistically significant, using a conservative statistical approach (see text).

Table IAntibody V_H and V_L (Kappa) Amino Acid Pairs with the Strongest Covariations (ϕ -values)

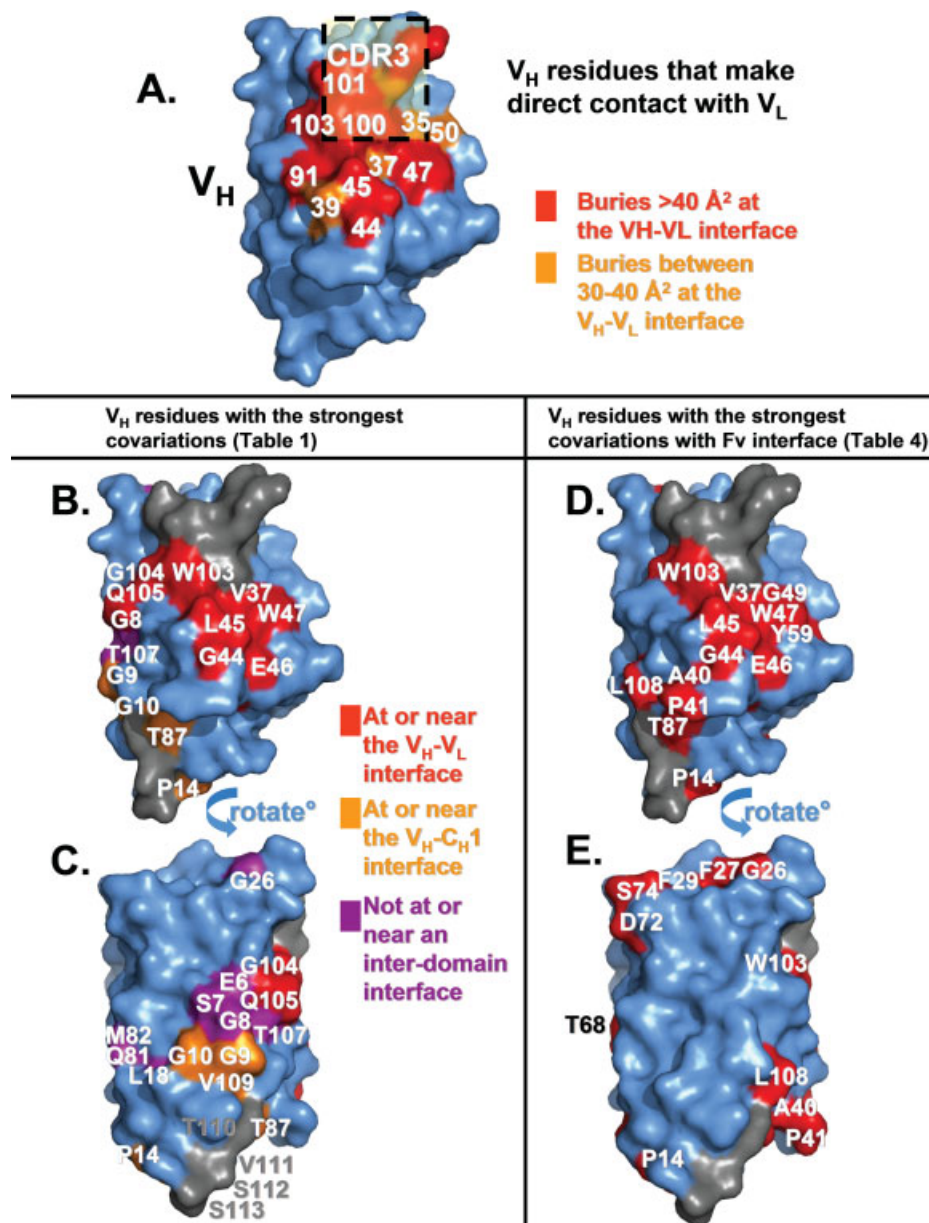
Top V_H covarying amino acids (A–B)	ϕ -value	V_H – V_L interface	V_H – C_H1 domain interface	Top V_L covarying amino acids (A–B)	ϕ -value	V_H – V_L interface	V_L – C_L domain interface
G9–L18	0.71		A	Q37–G64	0.64	A	
G9–G10	0.69		A–B	G57–G64	0.60		
L45–W47	0.68	A–B		F98–L104	0.54	A	B
E6–G9	0.65		B	P44–G64	0.53	A	
G8–G9	0.65		B	Q37–P59	0.52	A	
V37–W47	0.65	A–B		Y36–P44	0.48	A–B	
V63–M82	0.64			Q37–G57	0.48	A	
G104–G106	0.64	A		S63–G64	0.48		
S7–G8	0.62			Q37–S67	0.47	A	
V63–Q81	0.62			Q37–K39	0.46	A–B	
E6–L18	0.60			Q37–P44	0.46	A–B	
G26–W47	0.60	B		P44–P59	0.46	A	
G44–W47	0.60	A–B		G57–S67	0.46		
G44–L45	0.60	A–B		P59–S63	0.45		
E6–G8	0.59			Y36–L46	0.44	A–B	
G8–T87	0.59		B	Q37–G68	0.44	A	
Q81–M82	0.59			P44–I75	0.44	A	
W103–V109	0.59	A	B	P44–G57	0.43	A	
G8–L18	0.58			Q6–P8	0.42		
G106–T107	0.58			Y36–F98	0.42	A–B	
W103–Q105	0.58	A–B		Q37–I48	0.42	A–B	
G8–G26	0.57			Q37–S63	0.42	A	
G26–T87	0.57		B	I48–I75	0.42	A	
T87–W103	0.57	B	A	P44–S67	0.41	A	
G106–V109	0.57		B	G64–I75	0.41		
P14–W47	0.56	B	A	P8–Y36	0.40	B	
G26–E46	0.56	B		T5–Q37	0.40	B	
R38–E46	0.56	A–B		Q37–R54	0.40	A	
E46–W47	0.56	A–B		P59–I75	0.40		
E46–T87	0.56	A	B	P44–F98	0.39	A–B	
V37–L45	0.54	A–B		P44–S63	0.39	A	
G8–G10	0.53		B	I48–S63	0.39	A	

The two columns labeled “Top V_H [V_L] covarying amino acids” list the amino acids in the format A–B, and provide the residue codes and Kabat positions. Entries in the columns “ V_H – V_L Interface,” “ V_H – C_H1 domain interface,” and “ V_L – C_L domain interface” identify amino acids (A, B, or both of each pair) near the specified interface.

the covariations. After careful evaluation of the data, we designated ϕ -values between 0.25 and 0.5 as moderate covariations, and ϕ -values greater than 0.5 as strong covariations. Of the 4.1 million possible amino acid combinations within the V-class Ig-fold alignment, 3212 (0.078%) had ϕ -values ≥ 0.25 and 133 (0.003%) had ϕ -values ≥ 0.5 .

As validation of our covariation analysis, we examined the data in two ways to confirm that expected patterns were present in the results. First, we investigated the relationship between ϕ -values and distance between amino acid pairs in 3D space. Previous studies have shown that strongly covarying amino acid pairs often involve positions that are near each other in 3D structures, although the trend has invariably been reported as weak.^{23,40} To see if the same pattern was present in our data, we plotted the ϕ -values ≥ 0.3 against the distance between the amino acid pairs in 3D space using our Fab crystal structure. As reported by others, we found a weak but significant relationship between ϕ -value and the 3D proximity of the pair members (data not shown). Second, we investigated whether our covariation data recapitulated an

amino acid network known to exist within a particular subset of V-class Ig-folds. The example we chose was a set of five residues (residues 6–10, Kabat numbering⁴⁵) at the N-terminus of human/murine IgG V_H domains. These residues adopt different backbone conformations depending on the presence of specific amino acid pairs.⁴⁶ Four conformations exist, depending on whether glutamic acid or glutamine is present at V_H position 6. Q6 is not well-conserved within V_H domains (it is also found commonly in V_L domains) and does not have significant covariations. However, E6 is highly conserved in variable heavy chains (V_H2 , V_H3 , and V_H4 subclasses in particular). E6 correlations with residues 7–10 yield some of the highest ϕ -values of the covariation dataset (S7 = 0.51; G8 = 0.59; G9 = 0.65; and G10 = 0.53; Table I). These high correlations among residues 6–10 are consistent with the known involvement of these residues in determining the N-terminal β -strand conformation of IgG V_H domains. High ϕ -values were also found among other positions known to be structurally important including V_H subfamily-dependent core positions 18, 63, 67, and 82 that have been described previously.⁴⁷

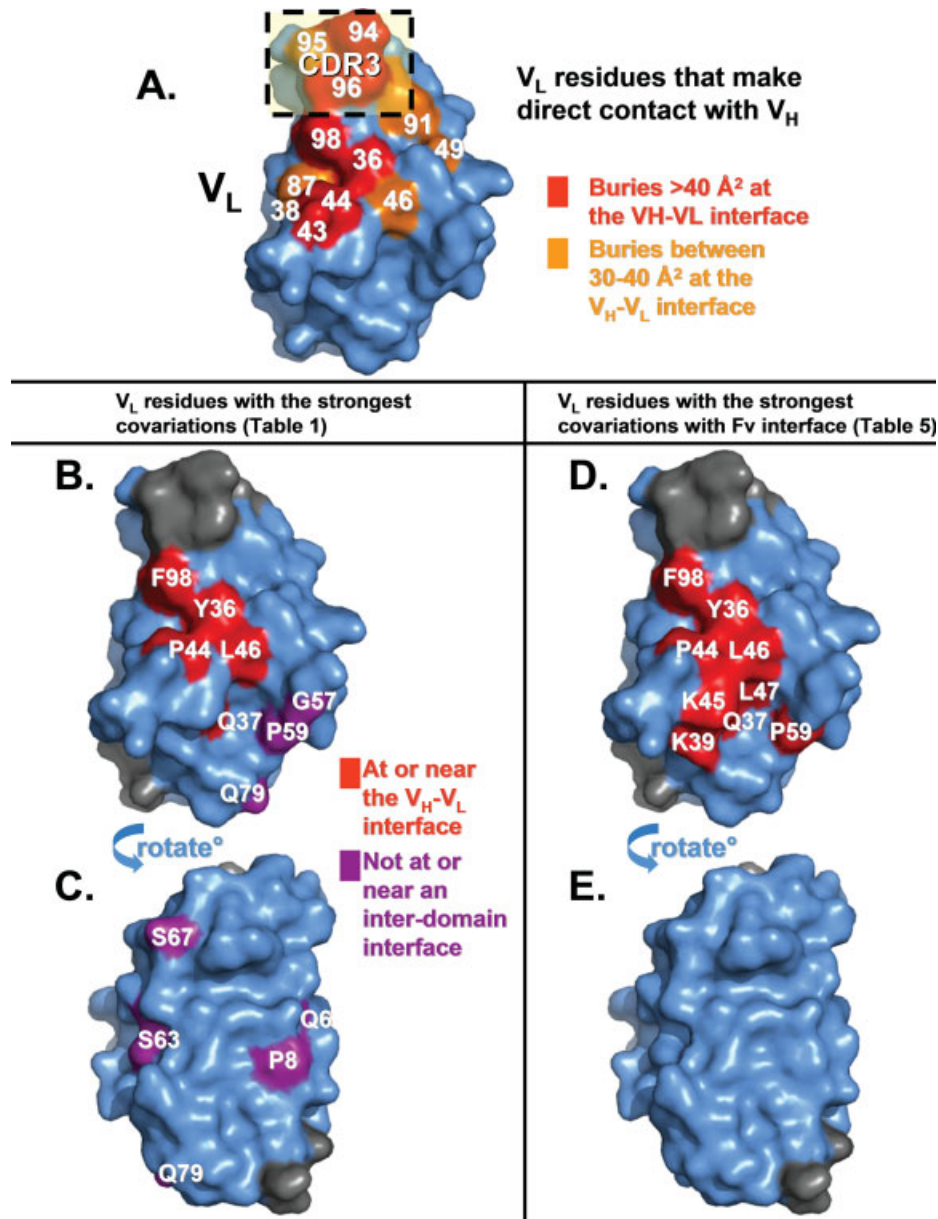
**Figure 3**

Covariations mapped to surface representations of an antibody V_H domain derived from an in-house Fab structure. (A) Surface representation of a V_H domain. Residues that bury $>40 \text{ \AA}^2$ at the Fv interface are shown in red and those that bury between 30 and 40 \AA^2 are shown in orange. In (B–E), residues colored grey (CDR3 residues as well as four residues at the C-terminus) were not match states in the HMM-derived V-class alignment and were not evaluated in this study. (B,C) V_H residues from amino acid pairs with the highest ϕ -values from Table I were mapped to the V_H surface: red, proximal to the V_H-V_L interface; orange, proximal to the V_H-C_H1 interface; and purple, distant from the two interfaces. (D,E) V_H residues from Table IV that display multiple covariations (ϕ -value >0.25) with V_H-V_L interface residues with greater than average ϕ -values are mapped onto the V_H surface in red. Residues from Tables I and IV that are completely buried in the interior of the structure are not shown.

Covariation results broadly applied to V_H and V_L domains

In this section, we describe patterns evident on examining the strongest covariations found within V_H and V_L domains. The V_H and V_L amino acid pairs with the highest ϕ -values are listed in Table I. The locations of the res-

idues involved in the strongest covariations from Table I were mapped onto our in-house Fab 3D structure [V_H , Fig. 3(B,C); V_L , Fig. 4(B,C)]. Interestingly, most of the residues contributing to the strongest covariations were found at or very near to the V_H-V_L interface [Table I, red in Figs. 3(B,C) and 4(B,C)], indicating a conserved amino acid network supporting this interface. V_H

**Figure 4**

Covariations mapped to surface representations of an antibody V_L domain derived from an in-house Fab structure. **A:** Surface representation of a V_L domain. Residues that bury $>40 \text{ \AA}^2$ at the Fv interface are shown in red and those that bury between 30 and 40 \AA^2 are shown in orange. In **(B-E)**, residues colored grey (CDR3 residues as well as four residues at the C-terminus) were not match states in the HMM-derived V-class alignment and were not evaluated in this study. **(B,C)** V_L residues from amino acid pairs with the highest ϕ -values from Table 1 were mapped to the V_L surface: red, proximal to the V_H - V_L interface; orange, proximal to the V_L - C_L interface; and purple, distant from the two interfaces. **(D,E)** V_L residues from Table V that display multiple covariations (ϕ -value >0.25) with V_H - V_L interface residues with greater than average ϕ -values are mapped onto the V_L surface in red. Residues from Tables I and V that are completely buried in the interior of the structure are not shown.

domains also appear to conserve an amino acid network near the variable-constant domain (V_H - C_H1) interface [Table I, orange in Fig. 3(B,C)]. This latter network, however, is not observed in V_L domains (Table I). Other strongly conserved residue pairs involved the N-terminal region of V_H (residues 6–10), and a few buried hydro-

phobic residues known to be highly subtype dependent (both described earlier).^{23,47}

For an alternative overview, we also compiled tables of V_H or V_L domain residues that had the most covariations (ϕ -values ≥ 0.25) with other amino acids in their respective domains (V_H , Table II; V_L , Table III), regardless of

Table II

Top 40 V_H Amino Acid Positions with the Most Covariations (ϕ -value > 0.25) with Other V_H Residues

Amino acid	Kabat#	#Links all	Avg. ϕ -value	#Interface links	Avg. ϕ -value to interface
G	10	74	0.39	3	0.34
G	8	74	0.35	4	0.37
T	87	71	0.38	6	0.44
W	103	69	0.36	6	0.36
M	82	69	0.39	1	0.41
G	26	67	0.37	6	0.46
Y	59	66	0.36	5	0.41
V	63	64	0.38	1	0.33
Q	81	63	0.36	2	0.31
W	47	61	0.37	6	0.54
E	46	61	0.36	6	0.41
R	19	60	0.36	1	0.33
L	18	60	0.38	3	0.36
I	69	56	0.34	5	0.32
T	68	56	0.33	4	0.38
G	49	56	0.33	5	0.43
E	6	56	0.37	1	0.39
I	51	55	0.35	5	0.38
Y	79	54	0.34	3	0.38
S	62	54	0.36	3	0.36
G	16	54	0.34	1	0.27
D	72	52	0.34	5	0.39
A	40	52	0.35	4	0.40
G	65	50	0.33	3	0.32
V	37	50	0.35	5	0.54
R	38	49	0.34	5	0.39
S	17	48	0.33	2	0.30
S	7	47	0.34	4	0.30
K	43	46	0.34	2	0.36
A	24	46	0.33	3	0.33
F	27	45	0.32	5	0.35
L	4	45	0.32	3	0.29
S	21	44	0.33	2	0.32
V	109	44	0.32	2	0.46
F	29	43	0.33	5	0.41
Q	105	42	0.32	3	0.38
R	71	39	0.32	1	0.28
S	25	39	0.32	3	0.38
L	82c	38	0.32	2	0.33
K	75	38	0.32	1	0.29

Residues that bury surface area at the Fv interface are highlighted in black rows. Residues immediately adjacent in primary sequence to interface residues are in grey rows.

the covariation strengths. Residues with the most covariations do not map to any single region of V_H (Supp. Info. Fig. 2) or V_L (Supp. Fig. 3), which is not surprising given the inclusion of many residues from weakly conserved networks. However, this analysis revealed an abundance of conserved networks evident with less stringent constraints on significant ϕ -values. We did note that residues involved in the most covariations mapped predominately to the β -sheet regions. Some covariations involve amino acids in the loop regions, but these are seen less frequently and likely reflect poorer alignment statistics in the loop regions, rather than a lack of conserved networks in the loop regions.

Analysis of the interface between antibody V_H and V_L domains

To further investigate V_H and V_L residues that may play a role in heavy and light chain association, we examined which V_H or V_L amino acids covary with amino acids that make direct contacts across chains at the V_H - V_L interface. Framework V_H and V_L positions that bury surface area at the V_H - V_L interface are in V_H 35, 37, 39, 44, 45, 47, 50, 91, and 103; and in V_L 36, 38, 43, 44, 46, 49, 87, and 98. These positions are mapped onto the surfaces of V_H and V_L in Figures 3(A) and 4(A), respectively, and onto a V-class sequence alignment (using our in-house V-class HMM) in Figure 5 (red letters in framework regions). The CDR3 loops of both V_H

Table III

Top 40 V_L Amino Acid Positions with the Most Covariations (ϕ -value > 0.25) with Other V_L Residues

Amino acid	Kabat#	#Links all	Avg. ϕ -value	#Interface links	Avg. ϕ -value to interface
G	64	47	0.4	6	0.36
G	57	44	0.34	6	0.33
S	22	44	0.33	0	0
V	104	44	0.32	0	0
P	59	43	0.35	6	0.35
Q	100	42	0.32	0	0
Q	37	41	0.37	6	0.35
P	44	40	0.33	6	0.35
S	65	36	0.31	0	0
S	67	36	0.35	5	0.33
G	68	35	0.32	4	0.28
I	75	35	0.32	4	0.36
I	48	34	0.33	4	0.35
Y	36	33	0.33	6	0.37
R	54	30	0.33	3	0.31
P	15	29	0.35	0	0
K	39	28	0.32	4	0.3
S	63	28	0.33	3	0.32
T	5	27	0.32	2	0.31
Q	79	25	0.31	3	0.3
P	8	23	0.31	4	0.32
Q	89	23	0.31	2	0.28
S	10	21	0.31	1	0.32
S	56	21	0.31	2	0.28
I	21	20	0.29	2	0.28
G	66	19	0.31	0	0
A	43	18	0.29	3	0.31
T	74	18	0.3	2	0.27
S	14	17	0.33	1	0.36
F	62	17	0.31	0	0
T	72	17	0.3	1	0.25
L	46	15	0.3	4	0.32
F	98	14	0.33	4	0.36
Q	6	13	0.3	1	0.39
Q	42	13	0.28	0	0
K	45	13	0.29	3	0.28
Y	49	13	0.29	2	0.27
V	58	11	0.3	2	0.25
D	70	11	0.3	0	0
A	84	11	0.31	0	0

Residues that bury surface area at the Fv interface are highlighted in black rows. Residues immediately adjacent in primary sequence to interface residues are in grey rows.

V_H region	Framework (FW) 1	CDR 1	Framework 2
Kabat # V_H	1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31	32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52	53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82
V_H 4 cons.	QLQESG--PGLVKPSETLSLTCTV-SG-GSI---SS---YYWSWIRQ-PP-G--K--GL--EWIGYIY		
V_H 2 cons.	QLKESG--PALVKPTQTTLTCTF-SG-FSL---STSG--VGVGWIRQ-PP-G--K--A1--EWLALID		
V_H 3 cons.	QLVESG--GGLVPPGGSRLRLSCAA-SG-FTF---SS---YAMSWVRQ-AP-G--K--GL--EWVSAIS		
V_H 5 cons.	QLVQSG--AEVKKPGESLKISCKG-SG-YST---TS---YWIGWVRQ-MP-G--K--GL--EWMGIY		
V_H 1a cons.	QLVQSG--AEVKKPGSSVKVSCA-SG-GTF---SS---YAIWVRQ-AP-G--Q--GL--EWMGGII		
V_H 1b cons.	QLVQSG--AEVKKPGASVKVSCA-SG-YTF---TS---YYMHWVRQ-AP-G--Q--GL--EWMGWIN		
V_H cons.	QLVESG--GGSVDAAGSLRLSCAA-SG-YTY---ST---YMGWVRQ-AP-G--K--GL--EWMGAIN		
V_H (HEWL)	QLQESG--PSLVKPSQTLSLTCSV-TG-SV---TS---YYWSWIRQ-PP-G--N--GL--EWMGYIS		
V_K 1 cons.	QMTQSPSSL-SASVGDRTVITCRA-S--QGI---S---SYLAWQQ-KP-G--K--A1--KLLIYAA		
V_K 2 cons.	VMTQSPSL-SLPVPGEFASISCRS-S--QSL--LHSN--GYNLWDW-LQ-KP-G--Q--S1--QLLIYLG		
V_K 3 cons.	VLTQSPATL-SLSPGERATLSCRA-S--QSV---SS---SYLAWQQ-KP-G--Q--A1--RLIYGA		
V_K 4 cons.	VMTQSPDSL-AVSLGERATINCRS-S--QSV--LYSS--NNKNYLAWQQ-KP-G--Q--P1--KLLIYWA		
V_{λ} 3 cons.	ELTQPP-SV-SVAPGQTARISCSG--DAL---GD---KYASWQQ-KP-G--Q--A1--VVIYDD		
V_{λ} 2 cons.	ALTQPA-SV-SGSPGQITISCTG-TS-SDV---GG---YNYVSWQQ-HP-G--K--A1--KMTYDV		
V_{λ} 1 cons.	VLTQPP-SV-SGAPGQRTISCSG-S-SNI---GS---NYVSWQQ-LP-G--T--A1--KLLIYDN		
Kabat # V_L	1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31	32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52	53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82
V_L region	Framework 1	CDR 1	Framework 2

V_H region	CDR 2	Framework 3	FW 4
Kabat # V_H	83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100 101 102 103 104 105 106 107 108 109 110 111 112	113 114 115 116 117 118 119 120 121 122 123 124 125 126 127 128 129 130 131 132 133 134 135 136 137 138 139 140 141 142 143 144 145 146 147 148 149 150 151 152 153 154 155 156 157 158 159 160 161 162 163 164 165 166 167 168 169 170 171 172 173 174 175 176 177 178 179 180 181 182 183 184 185 186 187 188 189 190 191 192 193 194 195 196 197 198 199 200 201 202 203 204 205 206 207 208 209 210 211 212 213 214 215 216 217 218 219 220 221 222 223 224 225 226 227 228 229 230 231 232 233 234 235 236 237 238 239 240 241 242 243 244 245 246 247 248 249 250 251 252 253 254 255 256 257 258 259 260 261 262 263 264 265 266 267 268 269 270 271 272 273 274 275 276 277 278 279 280 281 282 283 284 285 286 287 288 289 290 291 292 293 294 295 296 297 298 299 300 301 302 303 304 305 306 307 308 309 310 311 312 313 314 315 316 317 318 319 320 321 322 323 324 325 326 327 328 329 330 331 332 333 334 335 336 337 338 339 340 341 342 343 344 345 346 347 348 349 350 351 352 353 354 355 356 357 358 359 360 361 362 363 364 365 366 367 368 369 370 371 372 373 374 375 376 377 378 379 380 381 382 383 384 385 386 387 388 389 390 391 392 393 394 395 396 397 398 399 400 401 402 403 404 405 406 407 408 409 410 411 412 413 414 415 416 417 418 419 420 421 422 423 424 425 426 427 428 429 430 431 432 433 434 435 436 437 438 439 440 441 442 443 444 445 446 447 448 449 450 451 452 453 454 455 456 457 458 459 460 461 462 463 464 465 466 467 468 469 470 471 472 473 474 475 476 477 478 479 480 481 482 483 484 485 486 487 488 489 490 491 492 493 494 495 496 497 498 499 500 501 502 503 504 505 506 507 508 509 510 511 512 513 514 515 516 517 518 519 520 521 522 523 524 525 526 527 528 529 530 531 532 533 534 535 536 537 538 539 540 541 542 543 544 545 546 547 548 549 550 551 552 553 554 555 556 557 558 559 560 561 562 563 564 565 566 567 568 569 570 571 572 573 574 575 576 577 578 579 580 581 582 583 584 585 586 587 588 589 590 591 592 593 594 595 596 597 598 599 600 601 602 603 604 605 606 607 608 609 610 611 612 613 614 615 616 617 618 619 620 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 636 637 638 639 640 641 642 643 644 645 646 647 648 649 650 651 652 653 654 655 656 657 658 659 660 661 662 663 664 665 666 667 668 669 670 671 672 673 674 675 676 677 678 679 680 681 682 683 684 685 686 687 688 689 690 691 692 693 694 695 696 697 698 699 700 701 702 703 704 705 706 707 708 709 710 711 712 713 714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 754 755 756 757 758 759 760 761 762 763 764 765 766 767 768 769 770 771 772 773 774 775 776 777 778 779 780 781 782 783 784 785 786 787 788 789 790 791 792 793 794 795 796 797 798 799 800 801 802 803 804 805 806 807 808 809 810 811 812 813 814 815 816 817 818 819 820 821 822 823 824 825 826 827 828 829 830 831 832 833 834 835 836 837 838 839 840 841 842 843 844 845 846 847 848 849 850 851 852 853 854 855 856 857 858 859 860 861 862 863 864 865 866 867 868 869 870 871 872 873 874 875 876 877 878 879 880 881 882 883 884 885 886 887 888 889 890 891 892 893 894 895 896 897 898 899 900 901 902 903 904 905 906 907 908 909 910 911 912 913 914 915 916 917 918 919 920 921 922 923 924 925 926 927 928 929 930 931 932 933 934 935 936 937 938 939 940 941 942 943 944 945 946 947 948 949 950 951 952 953 954 955 956 957 958 959 960 961 962 963 964 965 966 967 968 969 970 971 972 973 974 975 976 977 978 979 980 981 982 983 984 985 986 987 988 989 990 991 992 993 994 995 996 997 998 999 1000 1001 1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 1015 1016 1017 1018 1019 1020 1021 1022 1023 1024 1025 1026 1027 1028 1029 1030 1031 1032 1033 1034 1035 1036 1037 1038 1039 1040 1041 1042 1043 1044 1045 1046 1047 1048 1049 1050 1051 1052 1053 1054 1055 1056 1057 1058 1059 1060 1061 1062 1063 1064 1065 1066 1067 1068 1069 1070 1071 1072 1073 1074 1075 1076 1077 1078 1079 1080 1081 1082 1083 1084 1085 1086 1087 1088 1089 1090 1091 1092 1093 1094 1095 1096 1097 1098 1099 1100 1101 1102 1103 1104 1105 1106 1107 1108 1109 1110 1111 1112 1113 1114 1115 1116 1117 1118 1119 1120 1121 1122 1123 1124 1125 1126 1127 1128 1129 1130 1131 1132 1133 1134 1135 1136 1137 1138 1139 1140 1141 1142 1143 1144 1145 1146 1147 1148 1149 1150 1151 1152 1153 1154 1155 1156 1157 1158 1159 1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177 1178 1179 1180 1181 1182 1183 1184 1185 1186 1187 1188 1189 1190 1191 1192 1193 1194 1195 1196 1197 1198 1199 1200 1201 1202 1203 1204 1205 1206 1207 1208 1209 1210 1211 1212 1213 1214 1215 1216 1217 1218 1219 1220 1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231 1232 1233 1234 1235 1236 1237 1238 1239 1240 1241 1242 1243 1244 1245 1246 1247 1248 1249 1250 1251 1252 1253 1254 1255 1256 1257 1258 1259 1260 1261 1262 1263 1264 1265 1266 1267 1268 1269 1270 1271 1272 1273 1274 1275 1276 1277 1278 1279 1280 1281 1282 1283 1284 1285 1286 1287 1288 1289 1290 1291 1292 1293 1294 1295 1296 1297 1298 1299 1300 1301 1302 1303 1304 1305 1306 1307 1308 1309 1310 1311 1312 1313 1314 1315 1316 1317 1318 1319 1320 1321 1322 1323 1324 1325 1326 1327 1328 1329 1330 1331 1332 1333 1334 1335 1336 1337 1338 1339 1340 1341 1342 1343 1344 1345 1346 1347 1348 1349 1350 1351 1352 1353 1354 1355 1356 1357 1358 1359 1360 1361 1362 1363 1364 1365 1366 1367 1368 1369 1370 1371 1372 1373 1374 1375 1376 1377 1378 1379 1380 1381 1382 1383 1384 1385 1386 1387 1388 1389 1390 1391 1392 1393 1394 1395 1396 1397 1398 1399 1400 1401 1402 1403 1404 1405 1406 1407 1408 1409 1410 1411 1412 1413 1414 1415 1416 1417 1418 1419 1420 1421 1422 1423 1424 1425 1426 1427 1428 1429 1430 1431 1432 1433 1434 1435 1436 1437 1438 1439 1440 1441 1442 1443 1444 1445 1446 1447 1448 1449 1450 1451 1452 1453 1454 1455 1456 1457 1458 1459 1460 1461 1462 1463 1464 1465 1466 1467 1468 1469 1470 1471 1472 1473 1474 1475 1476 1477 1478 1479 1480 1481 1482 1483 1484 1485 1486 1487 1488 1489 1490 1491 1492 1493 1494 1495 1496 1497 1498 1499 1500 1501 1502 1503 1504 1505 1506 1507 1508 1509 1510 1511 1512 1513 1514 1515 1516 1517 1518 1519 1520 1521 1522 1523 1524 1525 1526 1527 1528 1529 1530 1531 1532 1533 1534 1535 1536 1537 1538 1539 1540 1541 1542 1543 1544 1545 1546 1547 1548 1549 1550 1551 1552 1553 1554 1555 1556 1557 1558 1559 1560 1561 1562 1563 1564 1565 1566 1567 1568 1569 1570 1571 1572 1573 1574 1575 1576 1577 1578 1579 1580 1581 1582 1583 1584 1585 1586 1587 1588 1589 1590 1591 1592 1593 1594 1595 1596 1597 1598 1599 1600 1601 1602 1603 1604 1605 1606 1607 1608 1609 1610 1611 1612 1613 1614 1615 1616 1617 1618 1619 1620 1621 1622 1623 1624 1625 1626 1627 1628 1629 1630 1631 1632 1633 1634 1635 1636 1637 1638 1639 1640 1641 1642 1643 1644 1645 1646 1647 1648 1649 1650 1651 1652 1653 1654 1655 1656 1657 1658 1659 1660 1661 1662 1663 1664 1665 1666 1667 1668 1669 1670 1671 1672 1673 1674 1675 1676 1677 1678 1679 1680 1681 1682 1683 1684 1685 1686 1687 1688 1689 1690 1691 1692 1693 1694 1695 1696 1697 1698 1699 1700 1701 1702 1703 1704 1705 1706 1707 1708 1709 1710 1711 1712 1713 1714 1715 1716 1717 1718 1719 1720 1721 1722 1723 1724 1725 1726 1727 1728 1729 1730 1731 1732 1733 1734 1735 1736 1737 1738 1739 1740 1741 1742 1743 1744 1745 1746 1747 1748 1749 1750 1751 1752 1753 1754 1755 1756 1757 1758 1759 1760 1761 1762 1763 1764 1765 1766 1767 1768 1769 1770 1771 1772 1773 1774 1775 1776 1777 1778 1779 1780 1781 1782 1783 1784 1785 1786 1787 1788 1789 1790 1791 1792 1793 1794 1795 1796 1797 1798 1799 1800 1801 1802 1803 1804 1805 1806 1807 1808 1809 1810 1811 1812 1813 1814 1815 1816 1817 1818 1819 1820 1821 1822 1823 1824 1825 1826 1827 1828 1829 1830 1831 1832 1833 1834 1835 1836 1837 1838 1839 1840 1841 1842 1843 1844 1845 1846 1847 1848 1849 1850 1851 1852 1853 1854 1855 1856 1857 1858 1859 1860 1861 1862 1863 1864 1865 1866 1867 1868 1869 1870 1871 1872 1873 1874 1875 1876 1877 1878 1879 1880 1881 1882 1883 1884 1885 1886 1887 1888 1889 1890 1891 1892 1893 1894 1895 1896 1897 1898 1899 1900 1901 1902 1903 1904 1905 1906 1907 1908 1909 1910 1911 1912 1913 1914 1915 1916 1917 1918 1919 1920 1921 1922 1923 1924 1925 1926 1927 1928 1929 1930 1931 1932 1933 1934 1935 1936 1937 1938 1939 1940 1941 1942 1943 1944 1945 1946 1947 1948 1949 1950 1951 1952 1953 1954 1955 1956 1957 1958 1959 1960 1961 1962 1963 1964 1965 1966 1967 1968 1969 1970 1971 1972 1973 1974 1975 1976 1977 1978 1979 1980 1981 1982 1983 1984 1985 1986 1987 1988 1989 1990 1991 1992 1993 1994 1995 1996 1997 1998 1999 2000 2001 2002 2003 2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016 2017 2018 2019 2020 2021 2022 2023 2024 2025 2026 2027 2028 2029 2030 2031 2032 2033 2034 2035 2036 2037 2038 2039 2040 2041 2042 2043 2044 2045 2046 2047 2048 2049 2050 2051 2052 2053 2054 2055 2056 2057 2058 2059 2060 2061 2062 2063 2064 2065 2066 2067 2068 2069 2070 2071 2072 2073 2074 2075 2076 2077 2078 2079 2080 2081 2082 2083 2084 2085 2086 2087 2088 2089 2090 2091 2092 2093 2094 2095 2096 2097 2098 2099 2100 2101 2102 2103 2104 2105 2106 2107 2108 2109 2110 2111 2112 2113 2114 2115 2116 2117 2118 2119 2120 2121 2122 2123 2124 2125 2126 2127 2128 2129 2130 2131 2132 2133 2134 2135 2136 2137 2138 2139 2140 2141 2142 2143 2144 2145 2146 2147 2148 2149 2150 2151 2152 2153 2154 2155 2156 2157 2158 2159 2160 2161 2162 2163 2164 2165 2166 2167 2168 2169 2170 2171 2172 2173 2174 2175 2176 2177 2178 2179 2180 2181 2182 2183 2184 2185 2186 2187 2188 2189 2190 2191 2192 2193 2194 2195 2196 2197 2198 2199 2200 2201 2202 2203 2204 2205 2206 2207 2208 2209 2210 2211 2212 2213 2214 2215 2216 2217 2218 2219 2220 2221 2222 2223 2224 2225 2226 2227 2228 2229 2230 2231 2232 2233 2234 2235 2236 2237 2238 2239 2240 2241 2242 2243 2244 2245 2246 2247 2248 2249 2250 2251 2252 2253 2254 2255 2256 2257 2258 2259 2260 2261 2262 2263 2264 2265 2266 2267 2268 2269 2270 2271 2272 2273 2274 2275 2276 2277 2278 2279 2280 2281 2282 2283 2284 2285 2286 2287 2288 2289 2290 2291 2292 2293 2294 2295 2296 2297 2298 2299 2300 2301 2302 2303 2304 2305 2306 2307 2308 2309 2310 2311 2312 2313 2314 2315 2316 2317 2318 2319 2320 2321 2322 2323 2324 2325 2326 2327 2328 2329 2330 2331 2332 2333 2334 2335 2336 2337 2338 2339 2340 2341 2342 2343 2344 2345 2346 2347 2348 2349 2350 2351 2352 2353 2354 2355 2356 2357 2358 2359 2360 2361 2362 2363 2364 2365 2366 2367 2368 2369 2370 2371 2372 2373 2374 2375 2376 2377 2378 2379 2380	

Table IV V_H Residues with Multiple Covariations (ϕ -value > 0.25) with V_H Residues that Bury Surface Area at the Fv Interface

Amino acid	Kabat#	#Links w/ interface residues ^a	Avg. ϕ -value w/interface residues ^a	#Links w/all positions ^a	Avg. ϕ -value w/all positions ^a	$\Delta\phi$ -value (interface-all)	Interface
V(I)	37	5(3)	0.54 (0.30)	50 (19)	0.35 (0.34)	0.19	X
W	47	6	0.54	61	0.37	0.17	X
L	45	5	0.53	33	0.36	0.17	X
G	44	4	0.52	31	0.35	0.17	X
P	14	6	0.46	29	0.35	0.11	C _H 1
G	49	5	0.43	56	0.33	0.10	X - 1
G	26	6	0.46	67	0.37	0.09	
F	29	5	0.41	43	0.33	0.08	
S	74	5	0.39	26	0.32	0.07	
T	87	6	0.44	71	0.38	0.06	C _H 1
E	46	6	0.41	61	0.36	0.05	X - 1
Y	59	5	0.41	66	0.36	0.05	
R	38	5	0.39	49	0.34	0.05	X + 1
D	72	5	0.39	52	0.34	0.05	
A	40	4	0.40	52	0.35	0.05	X + 1
T	68	4	0.38	56	0.33	0.05	
P	41	6	0.33	32	0.30	0.03	
I	51	5	0.38	55	0.35	0.03	X + 1
F	27	5	0.35	45	0.32	0.03	
L	108	5	0.34	12	0.31	0.03	C _H 1
S	82b	4	0.32	24	0.29	0.03	
W	103	6	0.36	69	0.36	0.00	X
Nonspecific							
Y	32	5	0.30	15	0.29	0.01	
G	55	5	0.30	29	0.30	0.00	
I	69	5	0.32	56	0.34	-0.02	
S	7	4	0.30	47	0.34	-0.04	
G	8	4	0.35	74	0.39	-0.04	
L	11	4	0.30	30	0.31	-0.01	
L	11	4	0.30	30	0.31	-0.01	
Q	39	3	0.27	12	0.29	-0.02	X

Residues are sorted based on the difference between their average ϕ -value with interface residues versus their average ϕ -value with all positions within the V-class alignment. Residues that bury surface area at the interface are highlighted in black and marked with an X in the final column. Residues that are adjacent in primary sequence to interface residues are highlighted in dark grey and marked with an $X \pm 1$ in the final column. Residues at the C_H1 interface are in light grey rows. "Non-specific" residues (bottom of table) are those falling below an arbitrary cutoff above which residues appear to have strong, specific connections with interface residues. This cutoff was chosen based on a $\Delta\phi$ -value ([average ϕ -value with interface residues] - [average overall ϕ -value]) ≤ 0.01 . W103 was grouped with the specific interface residues because it is an interface residue.

^aMinimum ϕ -value cutoff of 0.25.

as do many of the residues adjacent in primary sequence. W47_{VH}, which incidentally is the V_H framework residue that buries the second highest amount of surface area at the interface, appears to be the central node of the V_H -side of the interface network based on the number and strength of its covariations with other interface residues, even though it is not at the center of the residues that make direct contact within the V_H - V_L interface [Fig. 3(D)]. Y36_{VL} and P44_{VL} appear to be the central nodes of the V_L -side of the interface network based on the same criteria [Fig. 4(D)]. P14_{VH}, T87_{VH}, and L108_{VH}, all near the V_H -C_H1 interface, also covary strongly with the V_H interface residues, suggesting that the maintenance of both the Fv and V_H -C_H1 interfaces are co-conserved traits. In contrast, covariations were weak for V_L residues in the proximity of the V_L -C_L interface.

Four V_H residues, V37, G44, L45, and W47, form a patch on the surface of V_H that interacts with V_L [Fig.

3(A-C)]. Each of these four residues has 30 or more ϕ -values >0.25 with other V_H residues; however, the ϕ -values between these four residues are collectively the strongest observed for each of these residues, with an average ϕ -value of 0.6 (see Tables I and IV). The side chains do not pack directly against one another or appear to interact strongly, suggesting that the residues covary for functional reasons—in this case enabling immunoglobulin heavy chain V_H domains to interact with immunoglobulin light chain V_L domains. In this context, the sidechains of W47, L45, and V37 together form a roughly flat hydrophobic surface that matches well with residues P44, F87, and F98 of V_L (Fig. 6). A fifth V_H residue—W103, one of the only other V_H residues burying surface area at the interface with V_L —also covaries with the four V_H residues, though with weaker ϕ -values averaging 0.38.

Two other V_H residues—R38 and E46—strongly covary with all five of the V_H residues discussed above, with

Table V*V_L* Residues with Multiple Covariations (ϕ -value > 0.25) with *V_L* Residues that Bury Surface Area at the Fv Interface

Amino acid	Kabat#	#Links w/ interface residues	Avg. ϕ -value w/interface residues ^a	#Links w/all positions	Avg. ϕ -value w/all positions ^a	$\Delta\phi$ -value (interface-all)	Interface
Y	36	6	0.37	33	0.33	0.04	X
I	75	4	0.36	35	0.32	0.04	
L	47	2	0.33	9	0.29	0.04	X + 1
F	98	4	0.36	14	0.33	0.03	X
P	44	6	0.35	40	0.33	0.02	X
L	46	4	0.32	15	0.30	0.02	X
P	59	6	0.35	43	0.35	0.00	
Q	37	6	0.36	41	0.37	−0.01	X + 1
I	48	4	0.32	34	0.33	−0.01	X + 1
K	39	4	0.29	28	0.32	−0.03	X + 1
K	45	3	0.28	13	0.39	−0.11	X − 1
Nonspecific							
G	57	6	0.33	44	0.34	−0.01	
S	67	5	0.33	36	0.35	−0.02	
K(R)	103	3 (1)	0.32 (0.27)	4 (1)	0.34 (0.27)	−0.02	
P	8	4	0.32	23	0.31	−0.01	
G	64	4	0.28	35	0.32	−0.04	
G	68	4	0.29	35	0.32	−0.03	
D	85	4	0.32	36	0.36	−0.04	
R	54	3	0.31	30	0.33	−0.02	
S	63	3	0.32	28	0.33	−0.01	
Q	79	3	0.30	25	0.31	−0.01	
S	56	2	0.28	21	0.31	−0.03	

Residues are sorted based on the difference between their average ϕ -value with interface residues versus their average ϕ -value with all positions within the V-class alignment. Residues that bury surface area at the interface are highlighted in black and marked with an X in the final column. Residues that are adjacent in primary sequence to interface residues are highlighted in grey and marked with an X \pm 1 in the final column. “Non-specific” residues (bottom of table) are those falling below an arbitrary cutoff above which residues appear to have strong, specific connections with interface residues. This cutoff was chosen based on a $\Delta\phi$ -value ([average ϕ -value with interface residues] − [average overall ϕ -value]) \leq 0.00. Q37, I48, K39, and K45 were grouped with the specific interface residues because they are adjacent in primary sequence to interface residues.

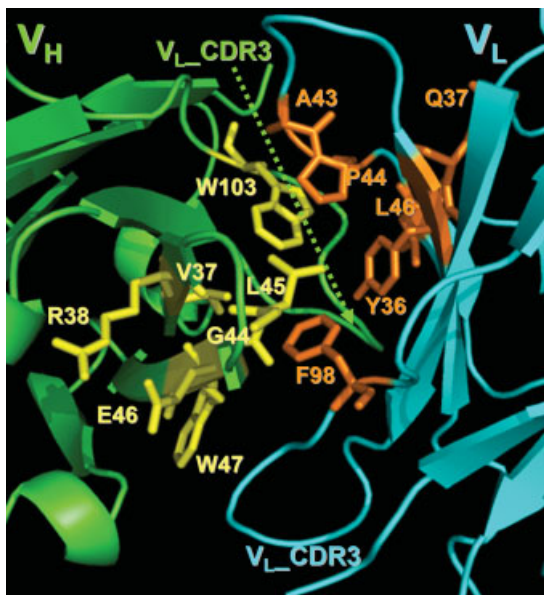
^aMinimum ϕ -value cutoff of 0.25.

average ϕ -values of 0.42 and 0.47, respectively. The ϕ -value between R38 and E46 is also strong [ϕ = 0.56, Table I, Fig. 3(A–C)]. R38 is almost completely buried by E46 in the interior of *V_H* and its guanidiny group forms a specific salt bridge with E46’s carboxyl group (Fig. 6). The charge burial of R38 is supported by other interactions with D90 and K/Q43. The E46 side chain does not contact any *V_H* residues at the *V_H*–*V_L* interface; thus, E46 is likely important for creating optimal surface topology and perhaps an electrostatic component important for *V_L* binding.

On the *V_L* side of the interface, the most strongly covarying interface residues are Y36 and P44 [the positional equivalents of *V_H* residues V37 and L45—Table I, Fig. 4(A–C)] with average ϕ -values with other *V_L* interface residues of 0.43 and 0.42, respectively. Also covarying with these two residues are Q37, A43, L46, and F98 (positional equivalents of *V_H* R38, G44, W47, and W103, Fig. 5), but with lower average ϕ -values (0.36, 0.31, 0.33, and 0.35, respectively). The ϕ -values between residues within the *V_L* domains were lower on average than those between *V_H* residues; this can be explained by the number of *V_L* sequences in the alignment being smaller (half

the number of the *V_H* sequences) and more heterogeneous (containing both *V_K* and *V_L* domains). Similar to what was observed for *V_H* interface residues, a cluster of *V_L* residues—Y36, A43, P44, L46 and F98—do not pack directly against one another, but instead combine to form the *V_H* binding surface. Unlike *V_H* residues R38 and E46 that form a salt bridge with one another, *V_L* positions 37 and 45 only weakly covary with one-another (Q37 and K45 have a ϕ -value = 0.34, Fig. 5). In general, the *V_L* residues important for *V_H*-binding are well conserved for both kappa and lambda light chain variable domains.

The strongly correlated residues at the Fv interface of both *V_H* and *V_L* thus appear to form conserved networks that enable recognition between the domains (Fig. 6). The interaction surface is fairly flat between the two domains. *V_H*–L45 and W103 insert themselves into a small groove created by *V_L* residues Y36, P44, L46, F87, and F98. The small size of *V_L*–P44 helps create the groove into which *V_H*–L45 and *V_H*–W103 intrude. In addition, the relatively small *V_H* V37 side chain creates a cavity on the hydrophobic surface of *V_H* into which the side chain of *V_L* F98 inserts (Fig. 6).

**Figure 6**

Structural view of the V_H and V_L residues from Tables IV and V. The polypeptide backbone of the V_H domain (green) and V_L domain (blue) are depicted using a cartoon ribbon diagram. V_H residues V37, R38, G44, L45, W47, and W103 are displayed in the stick format in yellow. V_L residues Y36, Q37, A43, P44, L46, and F98 are displayed in the stick format in orange.

Comparison of conserved residues networks of V_H domains and camelid V_{HH} domains lacking light chain interactions

Although most antibody V_H sequences associate with light chain V_L s, a subset of camelid variable heavy chain domains, denoted V_{HH} domains, are expressed naturally and function in the absence of both a light chain and a C_H1 domain.⁴⁹ V_{HH} domains are also substantially more soluble than V_H domains. The discovery of these simple and soluble V_H -like domains has had an enormous impact on antibody engineering because they represent potentially more facile reagents for protein design and discovery than traditional antibodies (which require combinations of heavy and light chains for function, stability, and solubility⁵⁰). We therefore investigated whether conserved residue networks differ between standard V_H and camelid V_{HH} domains. We expected to observe such differences at the positions in V_H that serve to support the V_H - V_L interface. Our results revealed 32 significant covariations involving identical positions within V_H and V_{HH} domains; however, the 32 covarying pairs contained different amino acids for V_H versus V_{HH} domains (Table VI). Among these contrasting residues are a tetrad of amino acids that have been previously reported to differentiate V_H from V_{HH} domains: V37F, G44E, L45R, and W47G.⁵¹ Substitution of this tetrad of camelid amino acids into V_H domains does not entirely impart

them with the solubility of V_{HH} domains; CDR3 composition and other factors have also been shown to be important.^{52–56} Our covariation results reveal additional framework residue positions, outside the tetrad described earlier, that naturally distinguish V_H from V_{HH} domains. These residues are at positions 13, 14, 33, 49, 63, 74, 82, 83, and 108. Solubilizing mutations at residues 74 and 108 have been reported.⁵⁷ Most of these positions are involved in networks surrounding the V_H - V_L or V_H - C_H1 interfaces (Table VI, Fig. 3), as expected. A consensus camelid V_{HH} sequence derived from the ~50 diverse camelid sequences in the V-class alignment is included in Figure 5 to highlight the positions of these observed differences between V_H and V_{HH} domains.

A natural human V_H raised against hen egg-white lysozyme was also demonstrated to be soluble in a monomeric form, similar to camelid domains (although the domain presumably maintains its ability to associate with V_L).^{52,58} This independently soluble anti-HEWL V_H domain contains yet another set of nonstandard V_H amino acids. Many of these residues are involved in the Fv interaction network and one is proximal to the V_H - C_H1

Table VI

Contrasting Features of V_H and Camelid V_{HH} Domains Based on Covariation Analyses

V_H linked pair	ϕ -value	Camelid V_{HH} linked pair	ϕ -value
G44–L45	0.61	E44–R45	0.57
L45–L108	0.29	R45–Q108	0.57
V(I)37–L45	0.54	F37–R45	0.50
P14–V37	0.40	A14–F37	0.50
L45–W47	0.70	R45–G47	0.46
V37–W47	0.65	F37–G47	0.44
—	—	C33 ^a –G47	0.44
—	—	A14–Q108	0.44
P14–G44	0.35	A14–E44	0.43
G44–W47	0.61	E44–G47	0.42
K13 ^b –L45	0.31	Q13–R45	0.42
V37–G44	0.53	F37–E44	0.40
—	—	C33 ^a –R45	0.40
L82 ^b –L45	0.29	M82–R45	0.38
V37–L108	0.32	F37–Q108	0.37
G49–L45	0.46	A49–R45	0.36
W47–L108	0.35	G47–Q108	0.36
L63 ^b –L45	0.30	V63–R45	0.36
—	—	C33 ^a –F37	0.36
P14–W47	0.48	A14–G47	0.35
—	—	C33 ^a –E44	0.35
K13 ^b –G44	0.30	Q13–E44	0.34
—	—	Q13–F37	0.32
S74–L45	0.45	A74–R45	0.30
L82 ^b –G44	0.29	M82–E44	0.27
S74–L108	0.25	A74–Q108	0.27
—	—	K83–E44	0.27
K13 ^b –W47	0.36	Q13–G47	0.26
—	—	V63–E44	0.26
S74–G44	0.38	A74–E44	0.25
S74–V37	0.33	—	—

^aC33 often makes a disulfide with CDR3 in camelid V_{HH} domains to stabilize camelid CDR3 structures.

^b V_{HH} 3 consensus matches the V_{HH} consensus amino acids at these positions.

interface: D27, D32, K39, K44, Y47, H59, K63, S68, and T108 (Fig. 5). Together with the V_{HH} results, it appears that multiple and independent amino acid networks may impart solubility to V_H and V_{HH} domains.

DISCUSSION

Despite an enormous amount of research involving antibodies and antibody-like therapeutics, very little use has been made of covariation analyses to investigate functional features of antibody domains. A study by Altshuh and coworkers^{59,60} investigated covariations across murine and human germline V_H or V_L sequences, with the goal of defining positions within each germline subclass that use mutually exclusive framework amino acid pairs to influence the structural conformations of CDR loops. Our study had a different goal: to use covariation analyses for determining naturally occurring amino acid networks, in antibody variable domains, that are generally important for antibody structure and function. To this end, we generated covariation data using a larger and more diverse set of V-class Ig-fold sequences.

Our results show that the most strongly conserved amino acid networks in antibody V_H and V_L domains are found at the interface between V_H and V_L , suggesting that preservation of this interface is a factor influencing antibody evolution. Interestingly, a small network of amino acids near the V_H - C_{H1} interface is also highly conserved. However, this network is not observed for residues near the V_L - C_L interface. Biophysical studies with light chains in isolation have shown that the V_L and C_L domains do not influence the unfolding kinetics or thermodynamics of one another, suggesting that the interaction between the two domains is weak.^{61,62} Alternately, Fabs (consisting of V_H , C_{H1} , V_L , and C_L domains) often show concerted unfolding reactions.⁴⁴ A viable explanation for the concerted unfolding reactions of Fabs compared with light chains in isolation is that the C_{H1} and C_L domains act as ideal linkers for the V_H and V_L domains and vice versa.⁶² It may also be that stronger interactions between V_H and C_{H1} , compared with V_L and C_L , promote cooperative unfolding of all four Fab domains.⁴⁴ The covariation data described in this report demonstrate that several V_H residues at the V_H - C_{H1} junction are involved in a strongly conserved network of amino acids and suggest that Variable-Constant domain interactions may be more important for immunoglobulin heavy chains than for immunoglobulin light chains.

CONCLUSIONS

In summary, we performed covariation analyses using a large, high-quality, and diverse alignment of V-class Ig-fold sequences. The data were used to discover antibody

variable domain amino acid networks that are evolutionarily conserved. Mapping the most highly conserved V_H and V_L networks to their structures revealed that the networks cluster near the V_H - V_L interface and near the V_H - C_{H1} interface, demonstrating the importance of preserving these interfaces during the evolution of antibody sequences.

These covariation data serve as a powerful tool for antibody (and Ig-fold domain) engineering. Insights from covariation analysis have improved our ability to rationally design more stable scFvs (Supp. Fig. 4). Most scFvs contain the majority of the residues within the V_H and V_L conserved networks revealed by the covariation data. Our initial approach for stabilizing scFvs has been to find gaps or holes within these existing networks that can be rectified by mutagenesis. Many of these changes have improved the thermal unfolding midpoint (T_M) of scFv V_H or V_L domains by 1–12°C (Miller *et al.*, manuscript in preparation). Stabilization of scFvs has enabled their use as building blocks that can be appended to full-length IgG molecules to produce stable bispecific IgG-like antibodies for consideration in clinical applications.⁶³

The covariation data described here may also be useful for engineering other aspects of V-class Ig-fold proteins, such as soluble V_H domains with conserved V_{HH} amino acid networks. The approach can be extended to other Ig-fold domains, such as C- and I-class, to identify amino acid networks supporting structure and function of other immunoglobulin or cell-surface receptor domains.

REFERENCES

1. Bird R, Hardmann K, Jacobson JW, Johnson S, Kaufman BM, Lee S, Lee T, Pope SH, Riordan GS, Whitlow M. Single-chain antigen-binding proteins. *Science* 1988;242:423–426.
2. Huston J, Levinson D, Mudgett-Hunter M, Tai M, Novotny J, Margolies MN, Ridge RJ, Brucoleri RE, Haber E, Crea R, Oppermann H. Protein engineering of antibody binding sites: recovery of specific activity of an anti-digoxin single-chain Fv analogue produced in *Escherichia coli*. *Proc Natl Acad Sci USA* 1988;85:5879–5883.
3. Glockshuber R, Malia M, Pfitzinger I, Plückthun A. A comparison of strategies to stabilize immunoglobulin Fv-fragments. *Biochemistry* 1990;29:1362–1367.
4. Brinkmann U, Reiter Y, Jung SH, Lee B, Pastan I. A recombinant immunotoxin containing a disulfide-stabilized Fv fragment. *Proc Natl Acad Sci USA* 1993;90:7538–7542.
5. Wörn A, Plückthun A. Stability engineering of antibody single-chain Fv fragments. *J Mol Biol* 2001;305:989–1010.
6. Demarest S, Glaser SM. Antibody therapeutics, antibody engineering, and the merits of protein stability. *Curr Opin Biotechnol* 2008; 11:675–687.
7. Bork P, Holm L, Sander C. The immunoglobulin fold. *J Mol Biol* 1994;242:309–320.
8. Williams A. The immunoglobulin superfamily—domains for cell surface recognition. *Annu Rev Immunol* 1988;8:381–405.
9. Lefranc M, Giudicelli V, Ginestoux C, Bodmer J, Muller W, Bontrop R, Lemaître M, Malik A, Barbie V, Chaume D. IMGT, the international ImMunoGeneTics database. *Nucleic Acids Res* 1999;27:209–212.

10. Carter P, Merchant AM. Engineering antibodies for imaging and therapy. *Curr Opin Biotechnol* 1997;8:449–454.
11. Ewert S, Honegger A, Plückthun A. Stability improvement of antibodies for extracellular and intracellular applications: CDR grafting to stable frameworks and structure-based framework engineering. *Methods* 2004;34:184–199.
12. Steipe B. Consensus-based engineering for protein stability: from intrabodies to thermostable enzymes. *Methods Enzymol* 2004;388:176–186.
13. Presta L. Engineering antibodies for therapy. *Curr Opin Biotechnol* 2002;74:237–256.
14. Davidson A. Multiple sequence alignment as a guideline for protein engineering strategies. *Methods Mol Biol* 2004;340:171–181.
15. Demarest S, Chen G, Kimmel BE, Gustafson D, Wu J, Salbato J, Poland J, Elia M, Tan X, Wong K, Short J, Hansen G. Engineering stability into *Escherichia coli* secreted Fabs leads to increased functional expression. *Protein Eng Des Sel* 2006;19:325–336.
16. Demarest S, Rogers J, Hansen G. Optimization of the antibody C_H3 domain by residue frequency analysis of IgG sequences. *J Mol Biol* 2004;335:41–48.
17. Altschuh D, Vernet T, Berti P, Moras D, Najai K. Coordinated amino acid changes in homologous protein families. *Protein Eng* 1988;2:193–199.
18. Altschuh D, Lesk AM, Bloomer AC, Klug A. Correlation of co-ordinated amino acid substitutions with function in viruses related to tobacco mosaic virus. *J Mol Biol* 1987;193:693–707.
19. Valencia A, Pazos F. Computational methods for the prediction of protein interactions. *Curr Opin Struct Biol* 2002;12:368–373.
20. Socolich M, Lockless SW, Russ WP, Lee H, Gardner KH, Ranganathan R. Evolutionary information for specifying a protein fold. *Science* 2005;437:512–518.
21. Süel G, Lockless SW, Wall MA, Ranganathan R. Evolutionary conserved networks of residues mediate allosteric communication in proteins. *Nature Struct Biol* 2003;10:59–69.
22. Magliery T, Regan L. Beyond consensus: statistical free energies reveal hidden interactions in the design of a TPR motif. *J Mol Biol* 2004;343:731–745.
23. Larson S, Di Nardo AA, Davidson AR. Analysis of covariation in an SH3 domain sequence alignment: applications in tertiary contact prediction and the design of compensating hydrophobic core substitutions. *J Mol Biol* 2000;303:433–446.
24. Pollock D, Taylor WR. Effectiveness of correlation analysis in identifying protein residues undergoing correlated evolution. *Protein Eng* 1997;10:647–657.
25. Govindarajan S, Ness JE, Kim S, Mundorff EC, Minshall J, Gustafson C. Systematic variation of amino acid substitutions for stringent assessment of pairwise covariation. *J Mol Biol* 2003;328:1061–1069.
26. Brenner S, Koehl P, Levitt M. The ASTRAL compendium for sequence and structure analysis. *Nucl Acids Res* 2000;28:254–256.
27. Chandonia J, Hon G, Walker NS, Lo Conte L, Koehl P, Levitt M, Brenner SE. The ASTRAL compendium in 2004. *Nucl Acids Res* 2004;32:D189–D192.
28. Murzin AG, Brenner SE, Hubbard T, Chothia C. SCOP: a structural classification of proteins database for the investigation of sequences and structures. *J Mol Biol* 1995;247:536–540.
29. Liu Y, Eisenberg D. 3D domain swapping: as domains continue to swap. *Protein Sci* 2002;11:1285–1299.
30. Lassmann T, Sonnhammer EL. Quality assessment of multiple alignment programs. *FEBS Lett* 2002;529:126–130.
31. Yang A, Honig B. An integrated approach to the analysis and modeling of protein sequences and structures. I. Protein structural alignment and a quantitative measure for protein structural distance. *J Mol Biol* 2000;301:665–678.
32. Jennings A, Edge CM, Sternberg MJ. An approach to improving multiple alignments of protein sequences using predicted secondary structure. *Protein Eng* 2001;14:227–231.
33. Eddy S. Profile hidden Markov models. *Bioinformatics* 1998;14:755–763.
34. Finn R, Mistry J, Schuster-Böckler B, Griffiths-Jones S, Hollich V, Lassmann T, Moxon S, Marshall M, Khanna A, Durbin R, Eddy SR, Sonnhammer EL, Bateman A. Pfam: clans, web tools and services. *Nucl Acids Res* 2006;34:D247–D251.
35. Henikoff S, Henikoff JG. Position-based sequence weights. *J Mol Biol* 1994;243:574–578.
36. Miller RJ. Simultaneous statistical inference. Springer, editor; 1981. 299p.
37. Wrabl J, Grishin NV. Gaps in structurally similar proteins: towards improvement of multiple sequence alignments. *Proteins: Struct Funct Genet* 2004;5471–87.
38. Koradi R, Billeter M, Wüthrich K. MOLMOL: a program for display and analysis of macromolecular structures. *J Mol Graph* 1996;14:51–55.
39. Gobel U, Sander C, Schneider R, Valencia A. Correlated mutations and residue contacts in proteins. *Proteins: Struct Funct Genet* 1994;18:309–317.
40. Chelvanayagam G, Eggenschwiler A, Knecht L, Gonnet GH, Benner SA. An analysis of simultaneous variation in protein structures. *Protein Eng* 1997;10:307–316.
41. Neher E. How frequent are correlated changes in families of protein sequences? *Proc Natl Acad Sci USA* 1994;91:98–102.
42. Jones S, Marin A, Thornton JM. Protein domain interfaces: characterization and comparison with oligomeric protein interfaces. *Protein Eng* 2000;13:77–82.
43. Goigou V, Cuisinier AM, Tonnelle C, Moinier M, Fougereau M, Fumoux F. Human immunoglobulin VH and VK repertoire revealed by in situ hybridization. *Mol Immunol* 1990;27:935–940.
44. Garber E, Demarest SJ. A broad range of Fab stabilities within a host of therapeutic IgGs. *Biochem Biophys Res Commun* 2007;355:751–757.
45. Wu T, Kabat EA. An analysis of the sequences of the variable regions of Bence Jones proteins and myeloma light chains and their implications for antibody complementarity. *J Exp Med* 1970;132:211–250.
46. Honegger A, Plückthun A. The influence of the buried glutamine or glutamate residue in position 6 on the structure of immunoglobulin variable domains. *J Mol Biol* 2001;309:687–699.
47. Honegger A. Engineering antibodies for stability and efficient folding. In: Chernajovsky Y, Nissim A, editors. *Therapeutic antibodies handbook of experimental pharmacology*. Volume 181. Berlin: Springer-Verlag; 2008 p. 47–68.
48. Honegger A, Plückthun A. Yet another numbering scheme for immunoglobulin variable domains: an automatic modeling and analysis tool. *J Mol Biol* 2001;309:657–670.
49. Hamers-Casterman C, Atarhouch T, Muyldermans S, Robinson G, Hamers C, Songa EB, Bendahman N, Hamers R. Naturally occurring antibodies devoid of light chains. *Nature* 1993;363:446–448.
50. Holliger P, Hudson PJ. Engineered antibody fragments and the rise of single domains. *Nature Biotechnol* 2005;23:1126–1136.
51. Reichmann L, Muyldermans S. Single domain antibodies: comparison of camel VH and camelised human VH domains. *J Immunol Methods* 1999;231:25–38.
52. Holt L, Herring C, Jespers LS, Woolven BP, Tomlinson IM. Domain antibodies: proteins for therapy. *Trends Biotechnol* 2003;21:484–490.
53. Desmyter A, Transue TR, Ghahroudi MA, Thi MH, Poortmans F, Hamers R, Muyldermans S, Wyns L. Crystal structure of a camel single-domain VH antibody fragment in complex with lysozyme. *Nature Struct Biol* 1996;3:803–811.
54. Decanniere K, Desmyter A, Lauwereys M, Ghahroudi MA, Muyldermans S, Wyns L. A single-domain antibody fragment in complex with RNase A: non-canonical loop structures and nanomolar affinity using two CDR loops. *Struct Fold Des* 1999;7:361–370.

55. Spinelli S, Frenken L, Bourgeois D, de Ron L, Bos W, Verrips T, Anguille C, Cambillau C, Tegoni M. The crystal structure of a llama heavy chain variable domain. *Nature Struct Biol* 1996;3:752–757.
56. Barthelemy P, Raab H, Appleton BA, Bond CJ, Wu P, Wiesmann C, Sidhu SS. Comprehensive analysis of the factors contributing to the stability and solubility of autonomous human VH domains. *J Biol Chem* 2008;283:3639–3654.
57. Tanha J, Nguyen T-D, Ng A, Ryan S, Ni F, MacKenzie R. Improving solubility and refolding efficiency of human V_Hs by a novel mutational approach. *Protein Eng Des Sel* 2006;19:503–509.
58. Li Y, Li H, Smith-Gill SJ, Mariuzza RA. Three-dimensional structures of the free and antigen-bound Fab from monoclonal antilysozyme antibody HyHEL-63. *Biochemistry* 2000;39:6296–6309.
59. Choulier L, Lafont V, Hugo N, Altschuh D. Covariance analysis of protein families: the case of the variable domains of antibodies. *Proteins: Struct Funct Genet* 2000;41:475–484.
60. Hugo N, LaFont V, Beukes M, Altschuh D. Functional aspects of co-variant surface charges in an antibody fragment. *Protein Sci* 2002;11:2697–2705.
61. Rowe E, Tanford C. Equilibrium and kinetics of the denaturation of a homogeneous human immunoglobulin light chain. *Biochemistry* 1973;12:4822–4827.
62. Röthlisberger D, Honegger A, Plückthun A. Domain interactions in the Fab fragment: a comparative evaluation of the single-chain Fv and Fab format engineered with variable domains of different stability. *J Mol Biol* 2005;347:773–789.
63. Glaser S, Demarest S, Miller BR, Wu X, Snyder WB, Wang N, Croner LJ; Biogen Idec MA Inc., assignee. Stabilized polypeptide compositions. PCT/US2008/0050370. 2007.